# An Empirical Study of Derivative-Free-Optimization Algorithms for Targeted Black-Box Attacks in Deep Neural Networks

**Giuseppe Ughi** · **Vinayak Abrol** · **Jared Tanner**

**Abstract** We perform a comprehensive study on the performance of derivative free optimization (DFO) algorithms for the generation of targeted black-box adversarial attacks on Deep Neural Network (DNN) classifiers assuming the perturbation energy is bounded by an $\ell_\infty$ constraint and the number of queries to the network is limited. This paper considers four pre-existing state-of-the-art DFO-based algorithms along with the introduction of a new algorithm built on BOBYQA, a model-based DFO method. We compare these algorithms in a variety of settings according to the fraction of images that they successfully misclassify given a maximum number of queries to the DNN. The experiments disclose how the likelihood of finding an adversarial example depends on both the algorithm used and the setting of the attack; algorithms limiting the search of adversarial example to the vertices of the $\ell^\infty$ constraint work particularly well without structural defenses, while the presented BOBYQA based algorithm works better for especially small perturbation energies. This variance in performance highlights the importance of new algorithms being compared to the state-of-the-art

G. Ughi
Mathematical Institute, University of Oxford
Andrew Wiles Building, Radcliffe Observatory Quarter
Oxford, OX2 6GG, United Kingdom
Tel.: +44 1865 273525
E-mail: ughi@maths.ox.ac.uk

V. Abrol
Mathematical Institute, University of Oxford
Andrew Wiles Building, Radcliffe Observatory Quarter
Oxford, OX2 6GG, United Kingdom
Tel.: +44 1865 273525
E-mail: abrol@maths.ox.ac.uk

J. Tanner
Mathematical Institute, University of Oxford
Andrew Wiles Building, Radcliffe Observatory Quarter
Oxford, OX2 6GG, United Kingdom
Tel.: +44 1865 273525
E-mail: tanner@maths.ox.ac.uk

in a variety of settings, and the effectiveness of adversarial defenses being tested using as wide a range of algorithms as possible.

**Keywords** Derivative Free Optimization · Deep Learning · Black-Box Attacks

## 1 Introduction

Deep Neural Networks (DNNs) achieve state-of-the-art performance on a growing number of applications such as acoustic modelling [22], image classification [20], and fake news detection [27] to name but a few. Alongside their growing application, there is a literature on the robustness of deep networks which shows that it is often possible to subtly perturb the input image of a DNN in order to degrade its performance; these perturbations are referred to as adversarial examples [15, 36]. For example, see [11, 14, 25, 35, 41] where road signals are perturbed so as to be wrongly interpreted by self driving cars that analyze images of them with DNNs. Methods to generate these adversarial examples are classified according to two main criteria [41]:

**Adversarial Specificity** establishes what the aim of the adversary is. In *non-targeted* attacks, the method perturbs the image in such a way that it is misclassified into any category other than the original one. While in *targeted* settings, the adversary specifies a category into which an image should be misclassified.

**Adversary's Knowledge** defines the amount of information available to the adversary. In *White-box* settings the adversary has complete knowledge of the network architecture and weights, while in the *Black-box* setting the adversary is only able to obtain the pre-classification output vector. The White-box setting allows for the use of gradients of a missclassification objective to efficiently compute the adversarial example [4, 8, 15], while the same optimization formulation of the Black-box setting requires use of a derivative free approach [2, 9, 23, 29].

In this work we consider the targeted black-box setting. In particular we follow [9] where:

- the *perturbation*, which causes the network to change the classification, is bounded in magnitude by a specified $\ell^\infty$-norm, $\varepsilon_\infty$, i.e. each pixel in the image cannot be perturbed by more than $\varepsilon_\infty$;
- the *number of queries* to the DNN needed to generate a targeted adversarial example should be as small as possible.

The Zeroth-Order-optimization (ZOO) algorithm proposed in [9] describes a Derivative Free optimization (DFO) method for computing adversarial examples in the black-box setting using a coordinate descent optimization method. At the time this was a substantial departure from previous black-box algorithms which trained a proxy DNN and then employ gradient based white-box attacks on the proxy network [32, 38]. It was demonstrated in [9] that these algorithms are especially effective when numerous adversarial examples are computed, but become less efficient when an individual adversarial examples is considered. Following the introduction of ZOO, there have been numerous improvements using other model-free DFO based approaches, see for example [1–3, 7, 23, 24, 28]. Many of these algorithms were developed in parallel, and so have not yet been bench-marked in a consistent setting, e.g. on the same network.

In this article, we present two frameworks for comparative evaluation of the existing algorithms that claim to have the fewest number of DNN queries to generate a successful
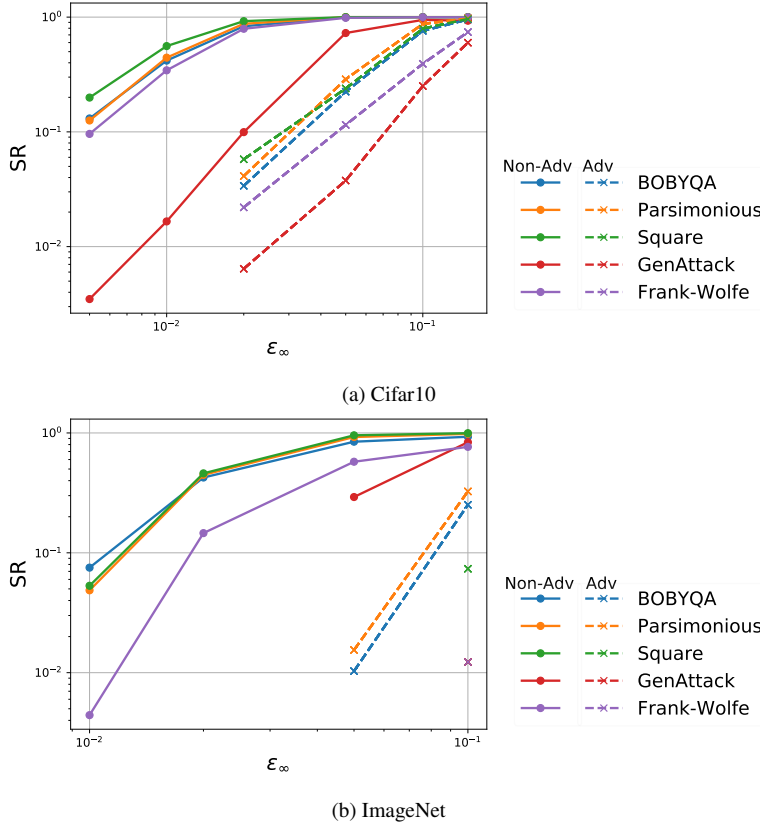
(a) Cifar10



(b) ImageNet

Fig. 1: The success rate (SR) of targeted attacks as a function of the perturbation's allowed $\ell^\infty$ magnitude for algorithms: GenAttack [2], Parsimonious [28], Square [3], Frank-Wolfe [7], and the BOBYQA based algorithm introduced here. Specifically for a ResNet50 network trained either on the CIFAR10 (a) or the ImageNet (b) dataset with (Adv) and without (Non-Adv) the defense by MadryLab [13]. An attack is considered successful if the method found the targeted adversarial example with less than 3'000 or 15'000 queries to the network trained on CIFAR and ImageNet dataset, respectively; Results for the case SR=0 i.e., when no perturbations were successful, are excluded from the plot.

attack. These are: GenAttack [2] which is based on a genetic direct-search method; Parsimonious algorithm [28], based on a combinatorial direct-search method on the vertices of the perturbation domain; the Square algorithm [3], based on a randomized direct-search method on the vertices of the perturbation domain; and the Frank-Wolfe algorithm [7] based on a momentum mechanism that approximates the gradient via finite differences. We also introduce a new algorithm built on a model-based DFO method [39]. In particular, we consider the *Bounded optimization BY Quadratic Approximation* (BOBYQA) [33] model-based DFO method which explicitly develops pseudo models to approximate the loss function in the optimization problem and then minimizes the loss function using methods from continuous optimization on the generated models. The aforementioned list of algorithms covers the

leading classes of DFO algorithms for limited function evaluations, see e.g., [10, 26] for recent reviews of DFO methods. The two frameworks are structured as follows:

1. In the first setting we consider attacks on DNNs trained on CIFAR10 and ImageNet datasets, with or without the adversarial defense by MadryLab [13]; this is the canonical setup for the comparison of black-box attacks that was considered in previous literature. We illustrate in Figure 1 a measure of how the performance of the considered algorithms compare, while further refined measures of comparison are included in Section 4. We observe that the algorithms that limit the optimization domain to the $\ell^\infty$ perturbation boundary, i.e. the Parsimonious and Square algorithms, are consistently the most effective. In particular, the Square algorithm achieves the highest Success Ratio (SR) with a fixed maximum number of queries, except for when the DNNs have been adversarially trained, and the Parsimonious algorithm achieves the highest SR when a network is trained with the MadryLab defense. However, these results are relative to the current state-of-the-art defense in a field which is in continuous development [12, 40] and newly proposed methods usually have a varying effect on the different attacking algorithms; for example the MadryLab defense [13] that we consider is most effective on Square algorithm in the ImageNet case.
2. In the second framework, the algorithms are allowed to perturb only a fraction of the pixels in the input; this is especially inspired by the structural defenses that transform the input in the wavelet space [18]. This framework allows us to understand the sensitivity of different algorithms to choices such as initialization, experimental protocol, dataset, and adversarial training. Our results demonstrate that the Parsimonious, Square, and BOBYQA based algorithms alternatively perform the best for different maximum perturbation energies.

The results in this paper show that the most likely algorithm to find an adversarial example varies according to the considered setting; the type of dataset, the defense, and the perturbation energy bound have a varying impact on the different algorithms. As a consequence of these experiments, new algorithms should be compared to the state-of-the-art in a variety of settings as done here, and the effectiveness of an adversarial defense should be tested with a variety of algorithms, including the BOBYQA based algorithm introduced in this paper.

The outline of the paper is as follows: in Section 2 we present how an adversarial example is generated by solving an optimization problem, and how DFO methods fit in this context. We also introduce the model-based BOBYQA algorithm. In Section 3 we present two popular techniques used in existing methods to improve the efficiency and scalability to high dimensional inputs. Section 4 presents the experimental setup and a comparative analysis of existing algorithms along with a focus on our proposed BOBYQA based algorithm. We close with some concluding remarks in Section 5.

## 2 Adversarial Examples Formulated as an optimization Problem

In classification tasks, a DNN outputs a vector whose length is equal to the number of classes and the DNN parameters are trained to match the maximum element of the given output to the correct class of the input. Adversarial perturbations are obtained by modifying the input in such a way that the maximum element of DNN output corresponds to a target class different from the original one.

Consider a classification operator $F : \mathscr{X} \to \mathscr{C}$ from input space $\mathscr{X}$ to output space $\mathscr{C}$ of classes. A targeted adversarial perturbation $\boldsymbol{\eta}$ to an input $\mathbf{X} \in \mathscr{X}$ has the property that it changes the classification to a specified target class $t$, i.e $F(\mathbf{X}) = c$ and $F(\mathbf{X} + \boldsymbol{\eta}) = t \neq c$.

Following the formulation in [2]; given an input space $\mathscr{X} = [l,u]^n$, with $l$ and $u$ being respectively the minimum and maximum values of the interval in which the pixels may vary, an output space $\mathscr{C} = \{1,\ldots,n_c\}$, where $n_c$ is the number of classes, a maximum energy budget $\varepsilon_\infty$, and a suitable loss function $\mathscr{L}$, then the task of computing the adversarial perturbation $\boldsymbol{\eta}$ can be cast as an optimization problem such as

$$\min_{\boldsymbol{\eta}} \mathscr{L}(\mathbf{X}, \boldsymbol{\eta}) \tag{1}$$

$$\text{s.t. } \|\boldsymbol{\eta}\|_\infty \leq \varepsilon_\infty;$$
$$[\mathbf{X} + \boldsymbol{\eta}]_j \geq l \qquad \forall j \in 1,\ldots,n$$
$$[\mathbf{X} + \boldsymbol{\eta}]_j \leq u \qquad \forall j \in 1,\ldots,n$$

where the final two inequality constraints are due to the perturbed image being still an image, i.e. $(\mathbf{X} + \boldsymbol{\eta}) \in \mathscr{X}$. Denoting the pre-classification output vector by $f(\mathbf{X})$, i.e. $F(\mathbf{X}) = \arg\max f(\mathbf{X})$, then the misclassification of $\mathbf{X}$ to target label $t$ is achieved by $\boldsymbol{\eta}$ if $f(\mathbf{X} + \boldsymbol{\eta})_t \geq \max_{j \neq t} f(\mathbf{X} + \boldsymbol{\eta})_j$. As demonstrated in [2, 4, 9], in this study we consider the following loss function for computing $\boldsymbol{\eta}$ in (1)

$$\mathscr{L}(\mathbf{X}, \boldsymbol{\eta}) = \log\left(\Sigma_{j \neq t} f(\mathbf{X} + \boldsymbol{\eta})_j\right) - \log\left(f(\mathbf{X} + \boldsymbol{\eta})_t\right). \tag{2}$$

Not having access to the internal parameters of the DNN, the gradient of the loss over the input space cannot be readily computed and instead the adversarial perturbation is found using specially adapted DFO algorithms.

## 2.1 Derivative Free optimization for Adversarial Examples

Derivative Free optimization is a well developed field with numerous classes of methods, see [10] and [26] for reviews on DFO principles and algorithms. Example classes of such methods include: direct search methods such as simplex, model-based methods, hybrid methods such as finite differences or implicit filtering, as well as randomized variants of the aforementioned and methods specific to convex or noisy objectives. For the generation of adversarial examples, the algorithms that we consider rely on three types of DFO methods:

- those where the gradient is computed via finite differences, either by sampling all the canonical directions as in ZOO attack [9] or random directions as in the Frank-Wolfe algorithm [7];
- those where the solution is thought to be in one of the vertices of the $\ell^\infty$ domain, i.e. $\boldsymbol{\eta}_i \in \{-\varepsilon_\infty, \varepsilon_\infty\}$ for any $i$. The Parsimonious algorithm [28] implements a combinatorial direct-search within the different possible vertices, initializing the perturbation to $-\varepsilon_\infty$ for all the pixels and then switching collections of them to $+\varepsilon_\infty$, when such an action decreases the loss function. The Square algorithm [3] instead implements a randomized direct-search method where square blocks of pixels are iteratively perturbed to be either $+\varepsilon_\infty$ or $-\varepsilon_\infty$;
- those where a direct search over the perturbation domain is performed using a genetic method such as GenAttack [2].

The optimization formulation in (1) is amenable to virtually all DFO methods, making it unclear which of the methods would be most effective in this context. Further, model-based methods are notably missing from the aforementioned list. Thus for completeness, we introduce an algorithm relying on a model-based method; specifically, BOBYQA is considered given its proven effectiveness in solving complex problems such as climate modelling [37].

2.2 Model-Based DFO

Given a set of $q$ samples $\mathscr{Y} = \{\mathbf{y}^1, ..., \mathbf{y}^q\}$ with $\mathbf{y}^i \in \mathbb{R}^n$, model-based DFO methods start by identifying the minimizer of the objective among the samples at iteration $k$, $\mathbf{x}^k = \arg\min_{\mathbf{y} \in \mathscr{Y}} \mathscr{L}(\mathbf{y})$. Following this, a model for the objective function $\mathscr{L}$ is constructed, typically centered around the minimizer. In its simplest form one uses a polynomial approximation to the objective, such as a quadratic model centered in $\mathbf{x}^k$

$$m_k(\mathbf{x}^k + \mathbf{p}) = a_k + \mathbf{c}_k^\top \mathbf{p} + \frac{1}{2}\mathbf{p}^\top \mathbf{M}_k \mathbf{p}, \tag{3}$$

with $a_k \in \mathbb{R}$, $\mathbf{c}_k, \mathbf{p} \in \mathbb{R}^n$, and $\mathbf{M}_k \in \mathbb{R}^{n \times n}$ being also symmetric. In a white-box setting one would set $\mathbf{c}_k = \nabla \mathscr{L}(\mathbf{x}^k)$ and $\mathbf{M}_k = \nabla^2 \mathscr{L}(\mathbf{x}^k)$, but this is not feasible in the black-box setting as we do not have access to the derivatives of the objective function. Thus at each iteration $k$, the parameters $a_k$, $\mathbf{c}_k$ and $\mathbf{M}_k$ are usually defined by imposing interpolation conditions

$$m_k(\mathbf{y}^i) = \mathscr{L}(\mathbf{y}^i) \quad \forall i \in 1, 2, \ldots, q, \tag{4}$$

and when $q < 1 + n + n(n+1)/2$ (i.e. the system of equations is under-determined) other conditions are introduced according to which method is considered. The objective model (3) is considered to be a good estimate of the objective in a neighborhood referred to as a trust region. Once the model $m_k$ is generated, the update step $\mathbf{p}$ is computed by solving the trust region problem

$$\min_{\mathbf{p}} \quad m_k(\mathbf{x}_k + \mathbf{p}) \tag{5}$$
$$\text{s.t.} \quad \|\mathbf{p}\| \leq \Delta,$$

where $\Delta$ is the radius of the region where we believe the model to be accurate, for more details see [31]. The new point $\mathbf{x}_k + \mathbf{p}$ is added to $\mathscr{Y}$ and a prior point is potentially removed. In this paper, we consider an exemplary model-based method called BOBYQA.

*2.2.1 BOBYQA*

The *Bound Optimization BY Quadratic Approximation* (BOBYQA) method, introduced in [33], updates the parameters of the model $a, \mathbf{c}$, and $\mathbf{M}$, in each iteration in such a way as to minimize the change in the quadratic term $\mathbf{M}_k$ between iterates while otherwise fitting the sample values:

$$\min_{a_k, \mathbf{c}_k, \mathbf{M}_k} \|\mathbf{M}_k - \mathbf{M}_{k-1}\|_F^2 \tag{6}$$
$$\text{s.t.} \quad m_k(\mathbf{y}^i) = \mathscr{L}(\mathbf{y}^i), \qquad \forall i \in 1, 2, \ldots, q,$$

with $n + 1 < q < 1 + n + n(n+1)/2$ and $\mathbf{M}_k$ initialized as the zero matrix. When the number of parameters $q = n + 1$ then the model is considered as linear with $\mathbf{M}_k$ set as zero. Every time a new query is done, the sample which is the least important geometrically is removed from $\mathscr{Y}$, thus keeping the dimension of $\mathscr{Y}$ fixed.

## 3 Improving Efficiency and Computational Scalability

Because of the high number of pixels in the input images, the generation of adversarial examples involves solving a high dimensional problem, which makes the use of any DFO method impractical; for instance, the application of the BOBYQA method requires the solution of (6) which scales in memory allocation at least quadratically with the input dimension, and thus is computationally too expensive. Consequently, the implementation of DFO based adversarial algorithms relies on strategies to reduce the dimensionality of the problem, this improves the computational scalability along with the efficiency, as demonstrated experimentally. Instead of solving (1) for $\boldsymbol{\eta} \in \mathbb{R}^n$ directly, the DFO based algorithms consider variations of the domain sub-sampling and/or hierarchical liftings techniques. Domain sub-sampling iteratively sweeps over batches of $b \ll n$ variables, while hierarchical lifting clusters and perturbs variables simultaneously, as described in following sections.

### 3.1 Domain Sub-Sampling

The simplest version of domain sub-sampling consists of partitioning the input dimension into smaller disjoint domains and optimizing the loss function in each of them sequentially. This is, in an $n$ dimensional problem, one considers $k = \lceil n/b \rceil$ sets of integers, $\{\Omega^j\}_{j=1}^k$, of size $b \ll n$ which are disjoint and which cover all of $[n]$. Then (1) is solved sequentially on the dimensions identified by the sets $\Omega^j$. This is possible since the optimization domain is box like, i.e. $\boldsymbol{\eta} \in [l,u]^n$, and each dimension's bound is independent from the others. Formally, rather than solving (1) for $\boldsymbol{\eta} \in \mathbb{R}^n$ directly, for each of $j = 1,\dots,k$ one *sequentially* solves for the $\boldsymbol{\eta}^j \in \mathbb{R}^n$ variables which are only non-zero for entries in $\Omega^j$. The resulting sub-domain perturbations $\boldsymbol{\eta}^j$ are then summed to generate the full perturbation $\boldsymbol{\eta} = \sum_{j=1}^k \boldsymbol{\eta}^j$, see Figure 2 as an example. That is, the optimization problem (1) is adapted to repeatedly looping over $j = 1,\dots,k$:

$$\min_{\boldsymbol{\eta}^j} \mathscr{L}\left(\mathbf{X} + \sum_{h \neq j} \boldsymbol{\eta}^\ell, \boldsymbol{\eta}^j\right) \tag{7}$$

$$\text{s.t.} \quad \left\| \sum_{h=1}^k \boldsymbol{\eta}^h \right\|_\infty \leq \varepsilon_\infty;$$

$$\left[ \mathbf{X} + \sum_{h=1}^k \boldsymbol{\eta}^h \right]_r \geq l \qquad \forall r \in \Omega^j;$$

$$\left[ \mathbf{X} + \sum_{h=1}^k \boldsymbol{\eta}^h \right]_r \leq u \qquad \forall r \in \Omega^j,$$

where the sets $\{\Omega^j\}_{j=1}^k$ are usually computed again once $j$ is equal to $k$, and the sub-domain perturbations $\boldsymbol{\eta}^j$ are initialized as null.

We identified three possible ways of selecting the sub-domains $\{\Omega^j\}_{j=1}^k$;

- In *Random Sampling* one considers at each iteration a different random sub-samplings of the domain, i.e. $k = 1$. The ZOO algorithm used this kind of sampling [9].
- In *Ordered Sampling* one generates a random disjoint partitioning of the domain, i.e. $k = \lceil n/b \rceil$ and $\Omega_j \cap \Omega_l = \emptyset$ for any $j$ and $l$. A new partitioning is generated when each variable has been optimized over once. This sampling is implemented in the Parsimonious algorithm.
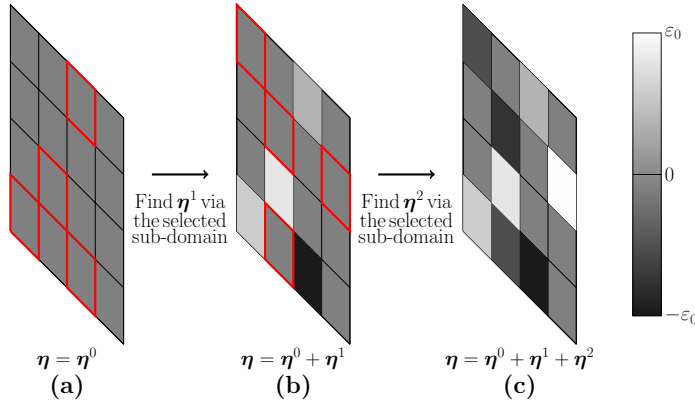
Fig. 2: Example of how the perturbation $\boldsymbol{\eta}$ evolves through the iterations when an image in $\mathbb{R}^{4\times4}$ is attacked. In (a) the perturbation is $\boldsymbol{\eta} = \boldsymbol{\eta}^0$ and a sub-domain of $b = 4$ pixels (in red) is selected. Once the optimal perturbation $\boldsymbol{\eta}^1$ in the selected sub-domain is found, the perturbation is updated in (b) and a new sub-domain of dimension $b$ is selected. The same is repeated in (c).

---

**Algorithm 1** GENERATE_SAMPLING_MATRIX($\hat{\mathbf{X}}$,$n_\ell$,b,j)

---

1: $\boldsymbol{\Omega} \leftarrow \mathbf{0} \in \mathbb{R}^{n_\ell \times b}$
2: $\mathbf{v} \leftarrow$ argsort Var($\hat{\mathbf{X}}$) # *Var defines the variance in intensity around a pixel.*
3: **for** $i = 1,\ldots,b$ **do**
4:     $\boldsymbol{\Omega}(\mathbf{v}[i + j \times b],[i]) = 1$.
5: **end for**
6: Return $\boldsymbol{\Omega}$.

---

- In *Variance Sampling* one still generates a a random disjoint partitioning of the domain, but chooses the sub-samplings sets $\{\Omega^j\}_{j=1}^k$ in order to optimize over the dimensions that have highest local variance in intensity first. Specifically, the variables are ordered by the variance in intensity among the 8 neighboring variables (e.g. pixels) in the same color channel of the input $\mathbf{X}$. The sets $\{\Omega^j\}_{j=1}^k$ are further reinitialized after each loop through $j = 1,\ldots,k$.

The sub-sampling of the domain affects the efficiency with which an algorithm successfully finds an adversarial example. For instance, in Figure 3 we compare how these different sub-sampling techniques affect the BOBYQA based algorithm when generating adversarial example for the MNIST and CIFAR10 dataset. It can be observed that variance sampling consistently has a higher success rate cumulative distribution function as compared with random and ordered sampling. This suggest that pixels belonging to high-contrast regions are more influential than the ones in low-contrast ones, and hence variance sampling is the preferable ordering.

To simplify the notation in the following section, the optimization variable is considered to be $\boldsymbol{\eta}^j = \boldsymbol{\Omega}^j \tilde{\boldsymbol{\eta}}^j$ where $\tilde{\boldsymbol{\eta}}^j \in \mathbb{R}^b$ and $\boldsymbol{\Omega}^j \in \mathbb{R}^{n \times b}$ is such that $[\boldsymbol{\Omega}^j]_{pq}$ is one if the $q$th element of $\Omega^j$ is $p$, zero otherwise. The implementation of variance sampling method at iteration $j$ in a domain of dimension $n_\ell$ is summarized in Algorithm 1.

(a) $\varepsilon_\infty = 0.4$

(b) $\varepsilon_\infty = 0.2$

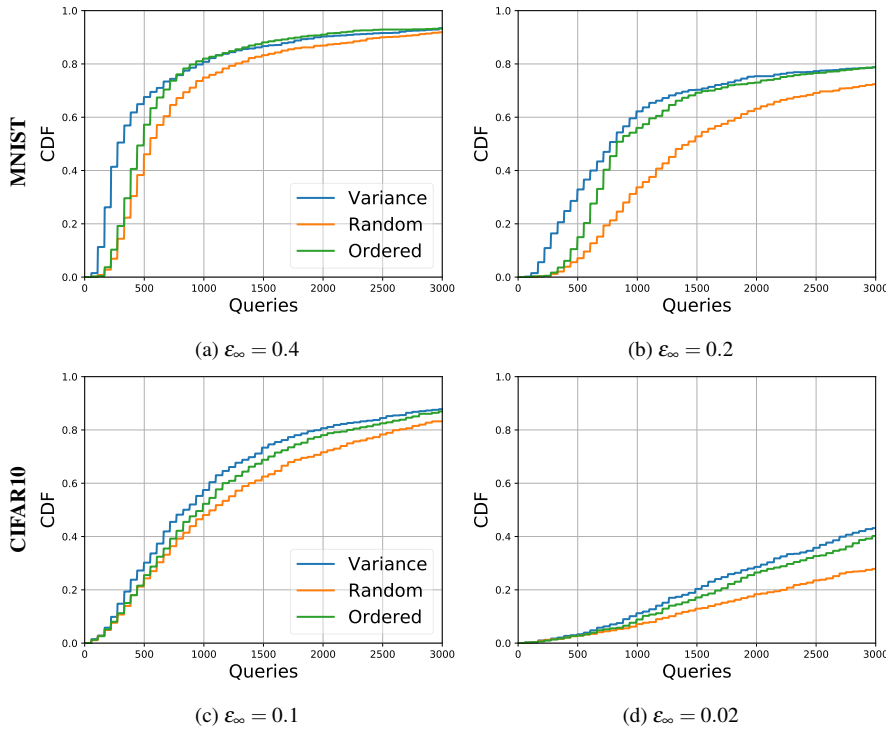(c) $\varepsilon_\infty = 0.1$

(d) $\varepsilon_\infty = 0.02$

Fig. 3: Cumulative distribution function of successfully perturbed images as a function of number of queries by the BOBYQA based algorithm attacking DNNs trained on the MNIST and the CIFAR10 datasets. In each image the effectiveness of different sub-sampling methods in generating a successful adversarial example is shown for different values of maximum perturbation energies $\varepsilon_\infty$. See [39] for details about experimental setup.

## 3.2 Hierarchical Lifting

Authors of ZOO attack [9] demonstrated that fewer queries are required to find adversarial example when pixels are considered in clusters, and not independently. This lead to the hierarchical lifting approach where one optimizes over increasingly higher dimensional spaces at each step, referred here as level $\ell$; Figure 4 shows how effective this approach is when implementing the BOBYQA based algorithm. These low dimensional spaces are lifted to the image space via a linear lifting, where at each level $\ell$ a linear lifting $\mathbf{D}^\ell : \mathbb{R}^{n_\ell} \to \mathbb{R}^n$ is considered and a perturbation $\hat{\boldsymbol{\eta}}_\ell \in \mathbb{R}^{n_\ell}$ is found to be added to the full perturbation $\boldsymbol{\eta}$, according to

$$\boldsymbol{\eta} = \sum_{j=0}^{\ell} \boldsymbol{\eta}_j = \sum_{j=0}^{\ell} \mathbf{D}^j \hat{\boldsymbol{\eta}}_j. \tag{8}$$

Here $\boldsymbol{\eta}_0$ is initialized as $\underline{\mathbf{0}}$ and the perturbations $\boldsymbol{\eta}_j$ of the previous layers are considered as fixed. An example of how this works is illustrated in Figure 5.

All the methods considered in this work rely on ideas which can be interpreted through this approach. The algorithms that we consider in this work rely on two kinds of linear lifting $\mathbf{D}^\ell$ differentiated by the way each scalar in $\hat{\boldsymbol{\eta}}$ is associated to a set of pixels in the original
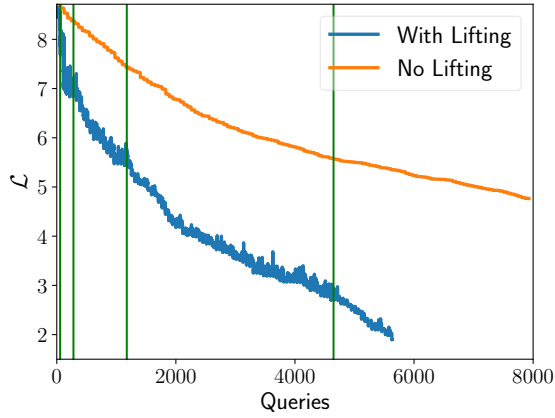
Fig. 4: Impact of hierarchical lifting approach on Loss function (2) as a function of the number of queries to Inception-v3 net trained on ImageNet dataset to find the adversarial example for a single image with the BOBYQA based method. The green vertical lines correspond to changes of hierarchical level, which entail an increase in the dimension of the optimization space.
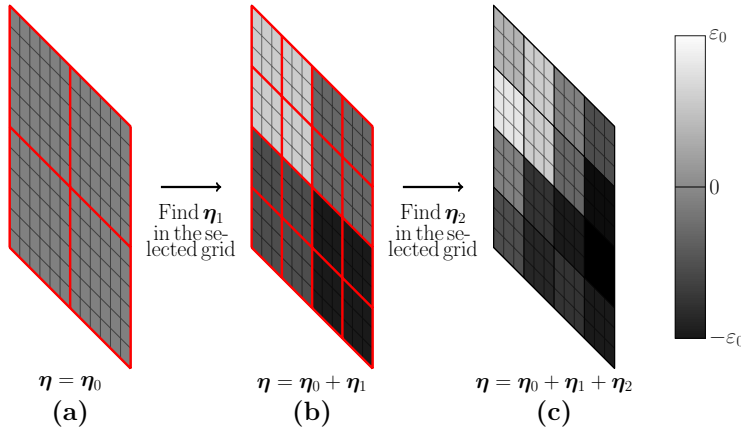


Fig. 5: Example of how the perturbation $\boldsymbol{\eta}$ is generated in a hierarchical lifting method with $n_1 = 4$ and $n_2 = 16$ on an image in $\mathbb{R}^{12 \times 12}$. In (a) the perturbation is $\boldsymbol{\eta} = \boldsymbol{\eta}_0$ and the boxes generated via the grid of dimension $n_1$ are highlighted in red. Once the optimal perturbation $\boldsymbol{\eta}_1$ is found, the perturbation is updated in (b) and the image is further divided with a grid with $n_2$ blocks. The final solution obtained after optimization is shown in (c).

image domain $\mathbb{R}^n$; namely the random and the block liftings. The former relates a random set of pixels of the original image to each hyper-variable; this forces the perturbation to be of high-frequency nature, as illustrated in Figure 6(a), which several articles indicate as being the most effective [16, 17, 34]. The GenAttack and Frank-Wolfe algorithms use a variation of this kind of lifting. The latter instead is based on interpolation operations; a sorting matrix $\mathbf{S}^\ell : \mathbb{R}^{n_\ell} \to \mathbb{R}^n$ is applied such that every index of $\hat{\boldsymbol{\eta}}_\ell$ is uniquely associated to a node of
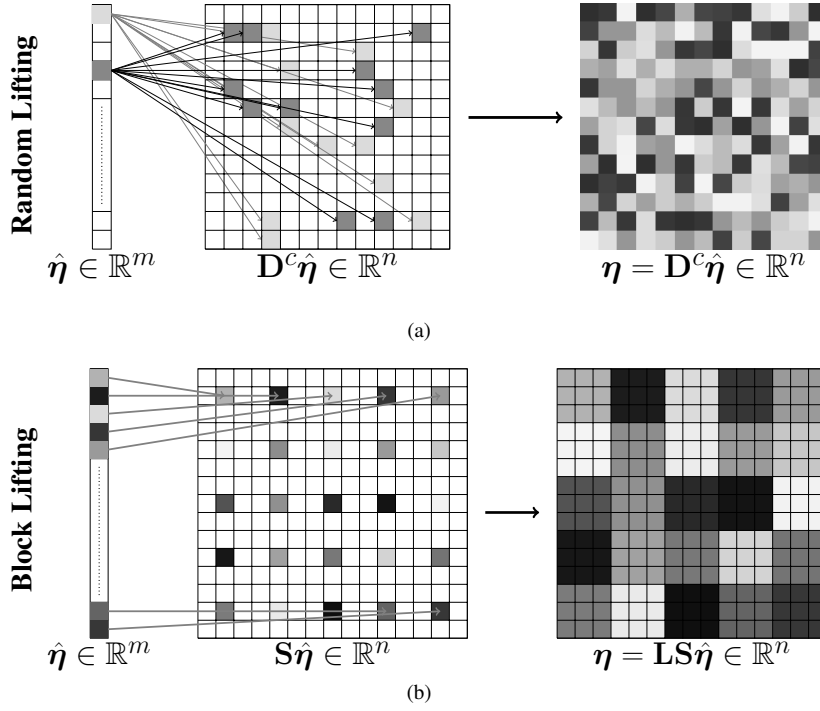
(a)



(b)

Fig. 6: Examples for (a) random and (b) block liftings. In the random case each pixel in the perturbation is associated to just one element of $\hat{\boldsymbol{\eta}}_\ell$. Block lifting uses a piece-wise constant interpolation $\mathbf{L}$ over a coarse grid $\mathbf{S}\hat{\boldsymbol{\eta}}_\ell$ and each block is associated uniquely to one of the variables in $\hat{\boldsymbol{\eta}}_\ell$. In both cases, the lifting $\mathbf{D}$ is such that each element $\mathbf{D}_{ij}$ is either 1 or 0.

a coarse grid masked over the original image. Afterwards, an interpolation $\mathbf{L}^\ell : \mathbb{R}^n \to \mathbb{R}^n$ is implemented over the values in the coarse grid, i.e. $\boldsymbol{\eta}_\ell = \mathbf{L}^\ell \mathbf{S}^\ell \hat{\boldsymbol{\eta}}_\ell = \mathbf{D}^\ell \hat{\boldsymbol{\eta}}_\ell$. Both Square and Parsimonious algorithms implement hierarchical lifting with the piece-wise constant interpolation, here referred to as block lifting. At the lower levels the interpolation lifting generates low frequency perturbations, as illustrated in Figure 6(b).

Since $n_\ell$ may still be very high, for each level $\ell$ domain sub-sampling is also applied considering $\hat{\boldsymbol{\eta}}_\ell = \sum_{j=0}^k \tilde{\boldsymbol{\eta}}_\ell^j$. In the piece-wise constant case with variance sampling, the blocks are ordered according to the variance of mean intensity among neighboring blocks, in contrast to the variance within each block as suggested in [9]. Consequently, at each level the adversarial example is found by solving the following iterative problem

$$\min_{\tilde{\boldsymbol{\eta}}_\ell^j} \quad \mathscr{L}\left(\mathbf{X}+\bar{\boldsymbol{\eta}}, \mathbf{D}^\ell \boldsymbol{\Omega}^k \tilde{\boldsymbol{\eta}}_\ell^j\right) \tag{9}$$
$$\text{s.t.} \quad \left\|\bar{\boldsymbol{\eta}}+\mathbf{D}^\ell \boldsymbol{\Omega}^k \tilde{\boldsymbol{\eta}}_\ell^j\right\|_\infty \le \varepsilon_\infty$$
$$\left[\mathbf{X}+\bar{\boldsymbol{\eta}}+\mathbf{D}^\ell \boldsymbol{\Omega}^k \tilde{\boldsymbol{\eta}}_\ell^j\right]_r \ge l \qquad \forall r \in \{1,...,n\}$$
$$\left[\mathbf{X}+\bar{\boldsymbol{\eta}}+\mathbf{D}^\ell \boldsymbol{\Omega}^k \tilde{\boldsymbol{\eta}}_\ell^j\right]_r \le u \qquad \forall r \in \{1,...,n\},$$

---

**Algorithm 2** GENERATE_LIFTING($n_\ell$,n)

---

1: $\mathbf{D} \leftarrow \mathbf{0} \in \mathbb{R}^{n \times n_\ell}$
2: **for** $i = 1, \ldots, n_\ell$ **do**
3:      Generate set of pixels $S$ that are in the block associated to the $i$-th element of the $n_\ell$ dimensional super-grid.
4:      **for** $j \in S$ **do**
5:          $\mathbf{D}(i, j) = 1$.
6:      **end for**
7: **end for**
8: Return $\mathbf{D}$.

---

where $\bar{\boldsymbol{\eta}} = \sum_{i=0}^{\ell-1} \boldsymbol{\eta}_i + \mathbf{D}^\ell \sum_{m \neq j} \hat{\boldsymbol{\eta}}_\ell^m$. Algorithm 2 gives an implementation of the block lifting matrix when in the grid has dimension $n_\ell$.

## 4 Comparison of Derivative Free Methods

In this section, we compare algorithms based on a selection of state-of-the-art DFO methods. In particular we consider BOBYQA based algorithm [39], GenAttack algorithm [2], Parsimonious algorithm [28], Square algorithm [3] and Frank-Wolfe algorithm [7] in the following two frameworks:

– Section 4.3 considers the canonical setup for black-box adversarial attacks on which the considered algorithms have been tuned in their respective articles. Specifically, we consider attacks on networks trained adversarially or not on CIFAR10 and ImageNet, two popular datasets in the literature, and with no further defense implemented.
– Section 4.4 considers a setup that simulates structural defenses on which the different algorithms were not tuned. We limit the perturbation to a fixed number of pixels with high variance in intensity considering attacks on a network non-adversarially trained on the CIFAR10 dataset.

The performance of all algorithms is measured in terms of the distribution of queries needed to successfully find adversaries to identical networks given a fixed $\ell^\infty$ perturbation constraint and the same input images.

### 4.1 Parameter Setup for Algorithms

The experiments use publicly available implementations for the GenAttack [2], Parsimonious [28], Square [3], and Frank-Wolfe [7] algorithms[1] using the same hyper-parameter setting and hierarchical lifting approach as suggested by the respective authors.

For the BOBYQA based algorithm [39], from Figure 3 we observed that the loss function is influenced the most by the pixels in high-contrast areas. Hence, we first apply the variance sub-sampling method followed by block lifting as described in Section 3.2[2]. Here, we

---

[1] GenAttack: `https://github.com/nesl/adversarial_genattack`
Parsimonious algorithm: `https://github.com/snu-mllab/parsimonious-blackbox-attack`
Square algorithm: `https://github.com/max-andr/square-attack`
Frank-Wolfe algorithm `https://github.com/uclaml/Frank-Wolfe-AdvML`

[2] The choice for this kind of lifting was driven by preliminary experiments in which we considered also a grid method with linear interpolation and a random lifting method as well. It is possible to run the analysis thanks to the code in [3]

---

**Algorithm 3** BOBYQA Based Algorithm

---

1: **Input:** Image $\mathbf{X} \in \mathbb{R}^n$, target label $t$, maximum perturbation $\varepsilon_\infty$, Neural Net $F$, initial hierarchical level grid dimensions $m$, maximum number of queries $n^{max}$, batch sampling size $b$, and maximum number $\kappa$ of queries that we are allowed to do for each batch.
2: **Initialize** $\boldsymbol{\eta} \leftarrow \underline{0} \in \mathbb{R}^n$, $n_{eval} = 0$, $\ell = 1$, $n_\ell = 12$.
3: **while** $\arg\max F(\mathbf{X} + \boldsymbol{\eta}) \neq t$ and $n_{eval} < n^{max}$ **do**
4:     *# Compute the number of sub samplings necessary to cover the whole domain*
5:     $num_{sub} = n/(n_\ell * b)$
6:     *# Generate the lifting matrix*
7:     $\mathbf{D}_\ell = \text{GENERATE\_LIFTING}(n_\ell, n)$
8:     *# Minimize on all the sampled sub-domains*
9:     **for** $j = 1, \dots, num_{sub}$ **do**
10:         *# Compute the matrix which selects $b$ dimensions of the $m$-dimensional domain.*
11:         $\boldsymbol{\Omega}_\ell^j = \text{GENERATE\_SAMPLING\_MATRIX}(\mathbf{X} + \boldsymbol{\eta}, n_\ell, b, j)$
12:         *# Define the pixel-wise bounds for a perturbation over $X + \boldsymbol{\eta}$.*
13:         $\mathbf{a} = \min\{l - \boldsymbol{\eta}, 0\}, \quad \mathbf{b} = \max\{u - \boldsymbol{\eta}, 0\}$
14:         *# Find $\hat{\boldsymbol{\eta}}_\ell^j$ by implementing the BOBYQA optimization to the problem (9).*
15:         $\hat{\boldsymbol{\eta}}_\ell^j = \text{BOBYQA}(F, \mathbf{X}, \boldsymbol{\eta}, \mathbf{a}, \mathbf{b}, \mathbf{D}_\ell, \boldsymbol{\Omega}_\ell^j, t)$     *# Algorithm 4*
16:         *# Update the noise*
17:         $\boldsymbol{\eta} + = \mathbf{D}_\ell \boldsymbol{\Omega}_\ell^j \hat{\boldsymbol{\eta}}_\ell^j$.
18:         $n_{eval} \mathrel{+}= \kappa$.
19:     **end for**
20:     $\ell + = 1$, $n_\ell * = 4$.
21: **end while**
22: **if** $\arg\max F(\mathbf{X} + \boldsymbol{\eta}) = t$ **then**
23:     The perturbation is successful.
24: **else if** $n_{eval} > n^{max}$ **then**
25:     The perturbation was not successful with $n^{max}$ iterations.
26: **end if**

---

---

**Algorithm 4** BOBYQA($F$, $\mathbf{X}$, $\boldsymbol{\eta}$, $\mathbf{a}$, $\mathbf{b}$, $\boldsymbol{\Omega}_\ell^j$, $\mathbf{D}^\ell$, $t$, $\kappa$)

---

1: Consider the restricted loss function $\mathscr{L}(\mathbf{X} + \boldsymbol{\eta}, \mathbf{D}^\ell \boldsymbol{\Omega}_\ell^j(\cdot)) : \mathbb{R}^b \to \mathbb{R}$
2: Build an initial model $m_0$ as in (3) of the loss function based on $b + 1$ samples; the samples consist of the initial perturbation $\mathbf{X} + \boldsymbol{\eta}$ and the $b$ perturbations obtained by considering changes along the canonical directions of $\mathbf{x}$ in $\mathbf{X} + \boldsymbol{\eta} + \mathbf{D}^\ell \boldsymbol{\Omega}_\ell^j \mathbf{x}$.
3: Find minimizer $\mathbf{x}$ of $m_o$ such that $\mathbf{D}^\ell \boldsymbol{\Omega}_\ell^j \mathbf{x} \in [\mathbf{a}, \mathbf{b}]$.
4: **for** $j = 1, \dots, \kappa - b$ **do**
5:     Add $\mathbf{x}$ to the set of samples and get rid of the least informative one according to [33].
6:     Build the new model $m_j$ according to (6).
7:     Find minimizer $\mathbf{x}$ of $m_j$ such that $\mathbf{D}^\ell \boldsymbol{\Omega}_\ell^j \mathbf{x} \in [\mathbf{a}, \mathbf{b}]$.
8: **end for**
9: Return $\mathbf{x}$.

---

consider an initial domain of dimension $n_1 = 2 \times 2 \times 3$, and double the refinement of the grid at each layer, i.e. $n_{\ell+1} = 4n_\ell$. Moreover, we observe for (6), the choice of a linear model to approximate the loss function works best, and we consequently consider the linear approximation in this paper; i.e., $\mathbf{M} = \mathbf{0}$ and $q = n + 1$ at all iterations, see [39]. The BOBYQA based algorithm is summarized in Algorithm 3 and a Python implementation of the proposed algorithm based on BOBYQA package from [6] is available on Github[3].

---

[3] https://github.com/giughi/An-Empirical-Study-of-DFO-Algorithms-for-Targeted-Black-Box-Attacks-in-DNNs

## 4.2 Dataset and Neural Network Specifications

We performed experiments using the popular ResNet50 architecture [21] with two training scenarios; one with the unperturbed images, and one with the defense[4] proposed in [13]. The number of experiments and the choice of the targets for each individual dataset is described below.

*CIFAR10* The CIFAR10 data-set contains images from 10 classes and of dimension 32x32x3. To generate a comprehensive distribution for the queries at each energy budget, ten correctly classified images are consider per each class, and each of them is targeted to all of the 9 remaining classes; this way we generate a total of 900 attacks per maximum perturbation energy per adversarial method.

*ImageNet* This data-set contains millions of images with a dimension of 299x299x3 divided among 1000 classes. Because of the high dimensionality and number of classes, random images are attacked considering a random target class. We conducted 200 and 160 tests for networks trained both with and without adversarial training per maximum perturbation energy.

## 4.3 Results for Standard and MadryLab Trained DNNs

In Figures 7 and 8 we present the cumulative fraction of images successfully misclassified (abridged by CDF for cumulative distribution function) as a function of the number of queries to the DNN for different maximum perturbation energies $\varepsilon_\infty$. The pixels are normalized to be in the interval $(-1/2, 1/2)$, hence, $\varepsilon_\infty = 0.1$ would imply that any pixel is allowed to change 10% of the total intensity range from its initial value. The CDFs are illustrated so that we can easily see which method has been able to misclassify the largest fraction of images in the given test-set for a fixed number of queries to the DNN.

For the CIFAR10 data-set in Figure 7, we observe that algorithms that search the perturbation directly in the vertices of the perturbation domain require the least amount of network queries. In the case of non-adversarially trained networks, the Square algorithm is able to misclassify using the least number of queries; this is demonstrated by its associated solid green CDF being consistently above that of the other methods. Specifically, when $\varepsilon_\infty = 0.05$, at 1,000 queries Square algorithms has a CDF of 0.97 compared to 0.94 and 0.88 of the Parsimonious and BOBYQA methods respectively, and for $\varepsilon_\infty = 0.005$ at 3,000 queries Square achieves a CDF of 0.20 which is 50% times higher than Parsimonious and BOBYQA. When the net is instead trained adversarially, dashed lines, Square algorithm looses a lot of its effectiveness becoming comparable to the BOBYQA based method, while Parismonious algorithm achieves almost always the highest fraction of successfully perturbed images for any given maximum number of queries. For example, when $\varepsilon_\infty = 0.05$ at 3,000 queries the CDF of Parisomonious is 0.29 compared to 0.25 and 0.23 of Square and BOBYQA.

In the ImageNet dataset, see Figure 8(a), we observe that an adversarial method can be especially susceptible to particular defenses. Specifically, when the network is trained without a defense, the Square algorithm has a success rate CDF that is consistently higher than the other methods, but the success rate CDF for the Square algorithm is decreased by the MadryLab defense so that it is substantially less effective than Parsimonious and BOBYQA

---

[4] These networks are available already trained at `https://github.com/MadryLab/robustness`
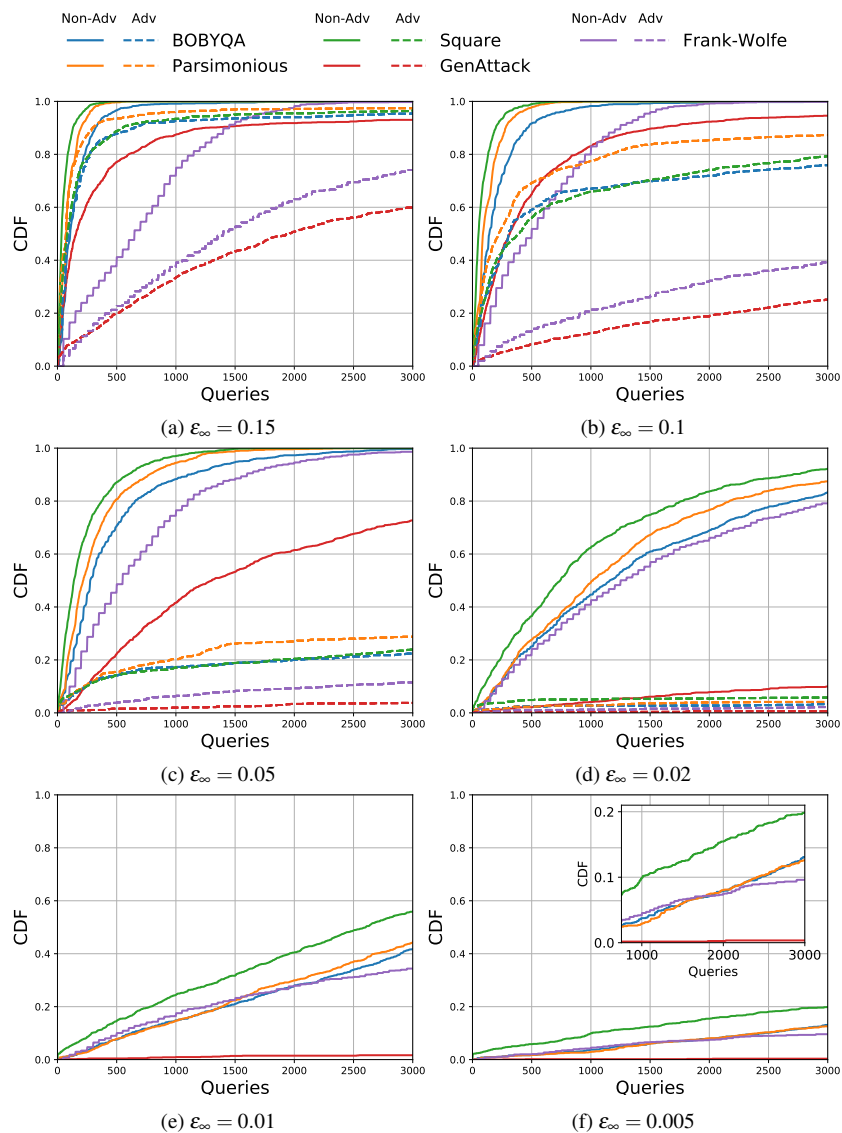
Fig. 7: Cumulative fraction of test set images successfully misclassified with adversarial examples generated by GenAttack, Parsimonious, Square, Frank-Wolfe, and our BOBYQA based approaches for different maximum perturbation energies $\varepsilon_\infty$ and DNNs trained on the CIFAR10 dataset. In all results the solid and dashed lines denoted by 'Non-Adv' and 'Adv' corresponds to attacks on networks trained without or with the MadryLab defense strategy [13] respectively.

algorithms. On the other hand, the Parsimonious method achieves similar results to Square algorithm in the non-adversarial case. On average for the different maximum perturbation energies Parsimonious is 0.045 less efficient than Square, but when the defense is introduced
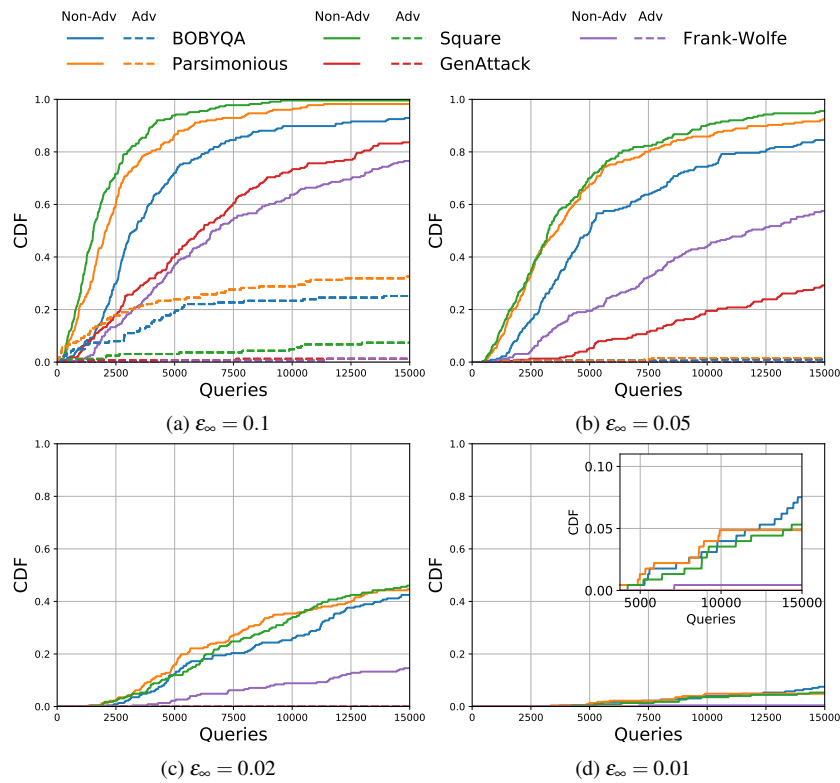
Fig. 8: Cumulative fraction of test set images successfully misclassified with adversarial examples generated by GenAttack, Parsimonious, Square, Frank-Wolfe, and our BOBYQA based approaches for different maximum perturbation energies $\varepsilon_\infty$ and DNNs trained on the ImageNet dataset. In all results the solid and dashed lines denoted by 'Non-Adv' and 'Adv' corresponds to attacks on networks trained without or with the MadryLab defense strategy [13] respectively.

it finds the adversarial examples with the least number of queries. In Figure 8(a) Parisomious has a CDF of 0.33 at 15,000 queries while BOBYQA 0.24 and Square 0.07. The rate with which the CDFs decrease as the maximum perturbation energy $\varepsilon_\infty$ decreases it also differs by algorithm. The CDF for Square decreases moderately faster than for Parsimonious such that Square has a consistently higher CDF than Parsimonious for $\varepsilon = 0.1$ in Figure 8(a) but consistently lower in Figure 8(d). Moreover, the success rate for BOBYQA decreases the slowest with $\varepsilon_\infty$ such that in Figure 8 its CDF is similar to or grater than Parsimonious. Specifically, in Figure 8(d) at 15,000 the final CDF of BOBYQA algorithm queries is 1.42 times higher than the one of the Square algorithm.

The Frank-Wolfe algorithm is able to achieve results comparable to the ones of the methods above while considering the small-dimensional problem of CIFAR10 with a very low maximum perturbation energy. However, when considering the ImageNet case and the adversarially trained DNNs, the Frank-Wolfe algorithm has a substantially lower success rate CDF; e.g. in the ImageNet case with non-adversarial training, Square algorithm achieves a CDF 1.66 times higher than the Frank-Wofle algorithm when $\varepsilon_\infty = 0.05$.

Finally, GenAttack has a higher success rate CDF than the Frank-Wolfe algorithm in the ImageNet case for $\varepsilon_\infty = 0.1$, see Figure 8(a), but, besides this case, it constantly achieves the lowest success rate.

## 4.4 Results with Fixed Pixel Count Constraints

In addition to network training designed to increase robustness, such as MadryLab considered previously, there are a multitude of other defenses and real world constraints [19]. The relative success rate, or other characteristics, of adversarial algorithms can be expected to differ in these diverse settings. To demonstrate this, we consider one such setting where the maximum number of pixels allowed to be perturbed is limited. This is motivated by the defenses where network inputs are thresholded in a wavelet domain to exclude high frequency perturbations [18], as well as by real world constraints such as attacks designed to appear structured such as localized perturbations designed to look like graffiti [14, 30]. We allow the algorithms to perturb only the fixed selection of the 1,000 pixels of the targeted image that have the highest variance in intensity in their channel neighborhood. Because of the previous results it is possible to identify three methods that work consistently better than the others, and thus only these will be considered, namely: the Parsimonious, the Square, and the BOBYQA based algorithms. To allow the perturbations to be limited to the selected pixels, we consider the Square algorithm with squares of pixel dimension, the Parsimonious algorithm on the finest grid, and the BOBYQA algorithm without the hierarchical lifting, i.e. $\mathbf{D}^1 = \mathbf{I}$ where $\mathbf{I}$ is the identity matrix.

The results reported in Figure 9 suggest that when the domain is dimensionally limited, the most efficient algorithm changes according to the allowed maximum perturbation energy. When the maximum perturbation energy decreases and the linear model is more accurate, the BOBYQA method manages to achieve a higher SR than both Square and Parsimonious algorithms, unlike in the previous experiments. Moreover, the Parsimonious algorithm has almost identical behavior to Square algorithm for high energy bounds, but becomes more efficient when the maximum energy is $\varepsilon_\infty = 0.05$. We also considered experiments on ImageNet, but limiting the number of pixels that could be perturbed did not allow for any successful misclassification with less than 15,000 queries.

## 5 Discussion and Conclusion

We have compared for the first time how the the existing GenAttack [2], Parsimonious [28], Square [3], and Frank-Wolfe [7] algorithms, and the newly introduced BOBYQA based method, behave when the available $\ell^\infty$ energy for a perturbation varies, and an adversarial training or a structural defense is considered.

The results suggest that those methods limiting the search for an adversarial example to the vertices of the $\ell^\infty$ perturbation domain generally work better. Whilst Square algorithm is especially effective on the non-adversarially trained networks, the Parsimonious algorithm manages to outperform any other approach when the networks are adversarially trained with the MadryLab implementation. Furthermore, the Parsimonious algorithm performs better than Square when considering the structural defense that limits the attacks on some pixels, suggesting that an algorithm based on combinatorial search is robust in its hyper-parameters to the setting where it is applied.
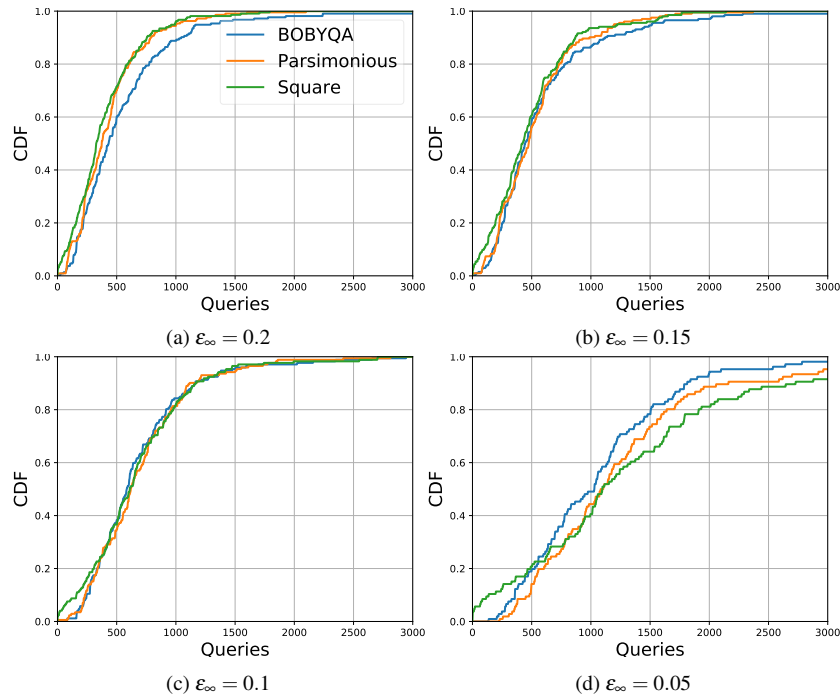
Fig. 9: Cumulative fraction of test set images successfully misclassified with adversarial examples generated by Parsimonious, Square, and our BOBYQA based approaches for different maximum perturbation energies $\varepsilon_\infty$ against a ResNet50 trained non-adversarially on the CIFAR10 dataset when only the 1000 pixels with the highest variance in intensity in their neighborhood are allowed to be modified.

The BOBYQA based algorithm was introduced in this paper to explore how model-based approaches compare to the state-of-the-art algorithms, and was found to achieve similar results to the Parsimonious and Square algorithms. In almost in all the experiments the BOBYQA based algorithm achieves a success rate CDF comparable to the ones of the Parsimonious and the Square algorithms; it achieves the state-of-the-art success rate at saturation for low maximum perturbation energy constraint both in the ImageNet case and in the pixel constrained problem. Moreover, new dimensionality reduction techniques that are being considered in DFO, see for example [5], might improve the results observed here and lead to a state-of-the-art algorithm for the generation of adversarial examples.

In conclusion, we find that both the structure of the algorithm and the attack setting have the potential to impact the algorithm performance. These observations highlight the importance of comparing any new algorithm to the state-of-the-art in a variety of different settings, such as is done here. Similarly, the effectiveness of an adversarial defense for DNNs should always be tested using as wide a range of algorithms as possible.

# References

1. Al-Dujaili, A., O'Reilly, U.M.: There are no bit parts for sign bits in black-box attacks. In: Proceedings of the International Conference on Learning Representations (ICLR) (2020)
2. Alzantot, M., Sharma, Y., Chakraborty, S., Zhang, H., Hsieh, C.J., Srivastava, M.B.: Genattack: Practical black-box attacks with gradient-free optimization. In: Proceedings of the Genetic and Evolutionary Computation Conference (GECCO), p. 1111–1119 (2019). DOI 10.1145/3321707.3321749
3. Andriushchenko, M., Croce, F., Flammarion, N., Hein, M.: Square attack: a query-efficient black-box adversarial attack via random search. Proceedings of the European Conference on Computer Vision (ECCV) (2020)
4. Carlini, N., Wagner, D.: Towards evaluating the robustness of neural networks. In: Proceedings of the IEEE Symposium on Security and Privacy (SP), pp. 39–57 (2017). DOI 10.1109/SP.2017.49
5. Cartis, C., Ferguson, T., Roberts, L.: Scalable derivative-free optimization for nonlinear least-square problems. Beyond First-Order Methods in ML Systems workshop at ICML (2020)
6. Cartis, C., Fiala, J., Marteau, B., Roberts, L.: Improving the flexibility and robustness of model-based derivative-free optimization solvers. ACM Trans. Math. Softw. **45**(3) (2019). DOI 10.1145/3338517
7. Chen, J., Zhou, D., Yi, J., Gu, Q.: A frank-wolfe framework for efficient and effective adversarial attacks. In: Proceedings fo the Association for the Advancement of Artificial Intelligence Conference (AAAI), pp. 3486–3494 (2020)
8. Chen, P.Y., Sharma, Y., Zhang, H., Yi, J., Hsieh, C.J.: Ead: elastic-net attacks to deep neural networks via adversarial examples. In: Proceedings of the AAAI Conference on Artificial Intelligence, pp. 10–17 (2018)
9. Chen, P.Y., Zhang, H., Sharma, Y., Yi, J., Hsieh, C.J.: ZOO: zeroth order optimization based black-box attacks to deep neural networks without training substitute models. In: Proceedings of the ACM Workshop on Artificial Intelligence and Security (AISec), p. 15–26 (2017). DOI 10.1145/3128572.3140448
10. Conn, A.R., Scheinberg, K., Vicente, L.N.: Introduction to derivative-free optimization, vol. 8. SIAM (2009)
11. Dalvi, N., Domingos, P., Sanghai, S., Verma, D., et al.: Adversarial classification. In: Proceedings of the ACM International conference on Knowledge Discovery and Data Mining (SIGKDD), pp. 99–108 (2004). DOI 10.1145/1014052.1014066
12. Dhillon, G.S., Azizzadenesheli, K., Lipton, Z.C., Bernstein, J., Kossaifi, J., Khanna, A., Anandkumar, A.: Stochastic activation pruning for robust adversarial defense. In: Proceedings of the International Conference on Learning Representations (ICLR) (2018)
13. Engstrom, L., Ilyas, A., Santurkar, S., Tsipras, D.: Robustness (python library) (2019). URL https://github.com/MadryLab/robustness
14. Eykholt, K., Evtimov, I., Fernandes, E., Li, B., Rahmati, A., Xiao, C., Prakash, A., Kohno, T., Song, D.: Robust physical-world attacks on deep learning visual classification. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pp. 1625–1634 (2018). DOI 10.1109/CVPR.2018.00175
15. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. In: Proceedings of the International Conference on Learning Representations (ICLR) (2015)

16. Gopalakrishnan, S., Marzi, Z., Madhow, U., Pedarsani, R.: Toward robust neural networks via sparsification. arXiv preprint arXiv:1810.10625 (2018)

17. Guo, C., Frank, J.S., Weinberger, K.Q.: Low frequency adversarial perturbation. In: Proceedings of the Conference on Uncertainty in Artificial Intelligence (UAI) (2018)

18. Guo, C., Rana, M., Cisse, M., van der Maaten, L.: Countering adversarial images using input transformations. In: Proceedings of the International Conference on Learning Representations (ICLR) (2018). URL https://openreview.net/forum?id=SyJ7ClWCb

19. Hao-Chen, H.X.Y.M., Deb, L.D., Anil, H.L.J.L.T., Jain, K.: Adversarial attacks and defenses in images, graphs and text: A review. International Journal of Automation and Computing **17**(2), 151–178 (2020)

20. He, K., Zhang, X., Ren, S., Sun, J.: Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In: Proceedings of the IEEE International Conference on Computer Vision (ICCV), p. 1026–1034 (2015). DOI 10.1109/ICCV.2015.123

21. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 770–778 (2016)

22. Hinton, G., Deng, L., Yu, D., Dahl, G.E., Mohamed, A., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Sainath, T.N., Kingsbury, B.: Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. IEEE Signal Processing Magazine **29**(6), 82–97 (2012). DOI 10.1109/MSP.2012.2205597

23. Ilyas, A., Engstrom, L., Athalye, A., Lin, J.: Black-box adversarial attacks with limited queries and information. In: Proceedings of the International Conference on Machine Learning (ICML), pp. 2137–2146 (2018)

24. Ilyas, A., Engstrom, L., Madry, A.: Prior convictions: Black-box adversarial attacks with bandits and priors. In: Proceedings of the International Conference on Learning Representations (ICLR) (2019)

25. Kurakin, A., Goodfellow, I., Bengio, S.: Adversarial examples in the physical world. In: Proceedings of the International Conference on Learning Representations (ICLR), Workshop Track (2017)

26. Larson, J., Menickelly, M., Wild, S.M.: Derivative-free optimization methods. Acta Numerica **28**, 287–404 (2019). DOI 10.1017/S0962492919000060

27. Monti, F., Frasca, F., Eynard, D., Mannion, D., Bronstein, M.M.: Fake news detection on social media using geometric deep learning. arXiv preprint arXiv:1902.06673 (2019)

28. Moon, S., An, G., Song, H.O.: Parsimonious black-box adversarial attacks via efficient combinatorial optimization. In: Proceedings of the International Conference on Machine Learning (ICML), pp. 4636–4645 (2019)

29. Narodytska, N., Kasiviswanathan, S.: Simple black-box adversarial attacks on deep neural networks. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops (CVPRW), pp. 1310–1318 (2017). DOI 10.1109/CVPRW.2017.172

30. Naseer, M., Khan, S., Porikli, F.: Local gradients smoothing: Defense against localized adversarial attacks. In: 2019 IEEE Winter Conference on Applications of Computer Vision (WACV), pp. 1300–1307 (2019)

31. Nocedal, J., Wright, S.J.: Numerical Optimization. Springer-Verlag New York (2006). DOI 10.1007/978-0-387-40065-5

32. Papernot, N., McDaniel, P., Goodfellow, I., Jha, S., Celik, Z.B., Swami, A.: Practical black-box attacks against machine learning. In: Proceedings of the ACM on Asia Conference on Computer and Communications Security (ASIA CCS), p. 506–519

(2017). DOI 10.1145/3052973.3053009

33. Powell, M.J.: The bobyqa algorithm for bound constrained optimization without derivatives. Tech. Rep. DAMTP 2009/NA06, University of Cambridge (2009)

34. Sharma, Y., Ding, G.W., Brubaker, M.: On the effectiveness of low frequency perturbations. Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI) pp. 3389–3396 (2019). DOI 10.24963/ijcai.2019/470

35. Sitawarin, C., Bhagoji, A.N., Mosenia, A., Chiang, M., Mittal, P.: Darts: Deceiving autonomous cars with toxic signs. arXiv preprint arXiv:1802.06430 (2018)

36. Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., Fergus, R.: Intriguing properties of neural networks. In: Proceedings of the International Conference on Learning Representations (ICLR) (2014)

37. Tett, S.F.B., Mineter, M.J., Cartis, C., Rowlands, D.J., Liu, P.: Can top-of-atmosphere radiation measurements constrain climate predictions? part i: Tuning. Journal of Climate **26**(23), 9348–9366 (2013). DOI 10.1175/JCLI-D-12-00595.1

38. Tu, C.C., Ting, P., Chen, P.Y., Liu, S., Zhang, H., Yi, J., Hsieh, C.J., Cheng, S.M.: Autozoom: Autoencoder-based zeroth order optimization method for attacking black-box neural networks. In: Proceedings of the AAAI Conference on Artificial Intelligence: Special Technical Track: AI for Social Impact (2019). DOI 10.1609/aaai.v33i01.3301742

39. Ughi, G., Abrol, V., Tanner, J.: A model-based derivative-free approach to black-box adversarial esxamples: Bobyqa. Proceedings of the Neural Information Processing Systems (NeruIPS) workshop "Beyond First Order Methods in ML" (2019)

40. Wang, X., Wang, S., Chen, P.Y., Wang, Y., Kulis, B., Lin, X., Chin, S.: Protecting neural networks with hierarchical random switching: Towards better robustness-accuracy trade-off for stochastic defenses. In: Proceedings of the International Joint Conference on Artificial Intelligence, IJCAI (2019)

41. Yuan, X., He, P., Zhu, Q., Li, X.: Adversarial examples: Attacks and defenses for deep learning. IEEE Transactions on Neural Networks and Learning Systems **30**(9), 2805–2824 (2019). DOI 10.1109/TNNLS.2018.2886017